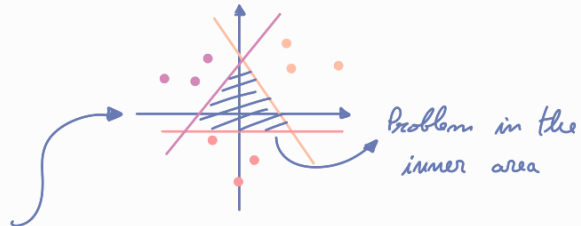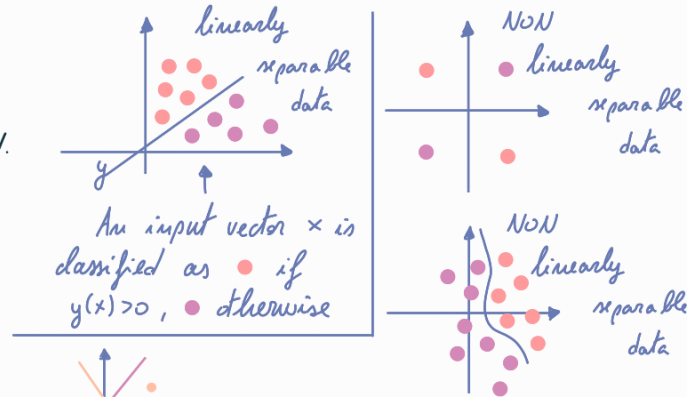# linearly separable data

instances in a dataset are linearly separable iff there exists a hyperplane that separates the instance space in two regions, such that differently classified instances are separated. All the samples of one class should be on one side. If the dataset is linearly separable there will exist infinite lines that can partition it.



linearly separable data

NON linearly separable data

NON linearly separable data

An input vector $x$ is classified as $\bullet$ if $y(x) > 0$, $\bullet$ otherwise

The solution of the problem will be:

- $y(x) = w^T x + w_0$     for two classes

- $y_i(x) = w_i^T x + w_{0i}$     for $K$ classes (we need $\forall i \in [1, K]$    to consider one linear model for each class)



Problem in the inner area

N.B. this is similar to linear regression but we are not giving a probabilistic interpretation of the solution; $y(x)$ is not anymore the prediction of the posterior probability of one class but it directly estimate the classification function

- $y(x) = w^T x + w_0 = \tilde{w}^T \tilde{x}$     with    $\tilde{w} = \begin{pmatrix} w_0 \\ w \end{pmatrix}$, $\tilde{x} = \begin{pmatrix} 1 \\ x \end{pmatrix}$     $\tilde{W}$ contains all the parameters of the model

- $y(x) = \begin{pmatrix} y_1(x) \\ \vdots \\ y_K(x) \end{pmatrix} = \begin{pmatrix} \tilde{w}_1^T \\ \vdots \\ \tilde{w}_K^T \end{pmatrix} \tilde{x} = \tilde{W}^T \tilde{x}$   with   $\tilde{W}^T = \begin{pmatrix} \tilde{w}_1^T \\ \vdots \\ \tilde{w}_K^T \end{pmatrix} = (\tilde{w}_1, \ldots, \tilde{w}_K)$

     directly the prediction of the output function (not prediction of the posterior)

This is called **linear model** because of course it is a linear combination. We are interested in computing $\tilde{W}$ (linear model in $\tilde{W}$).
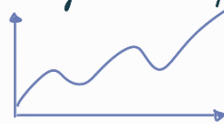
$$ y(\tilde{x}, \tilde{W}) = \tilde{W}^T \phi(x) $$

We can apply a transformation on the input such that the model will no longer be linear in $\tilde{x}$ but it will still be linear in $\tilde{W}$. We require linearity with respect to $\tilde{W}$.
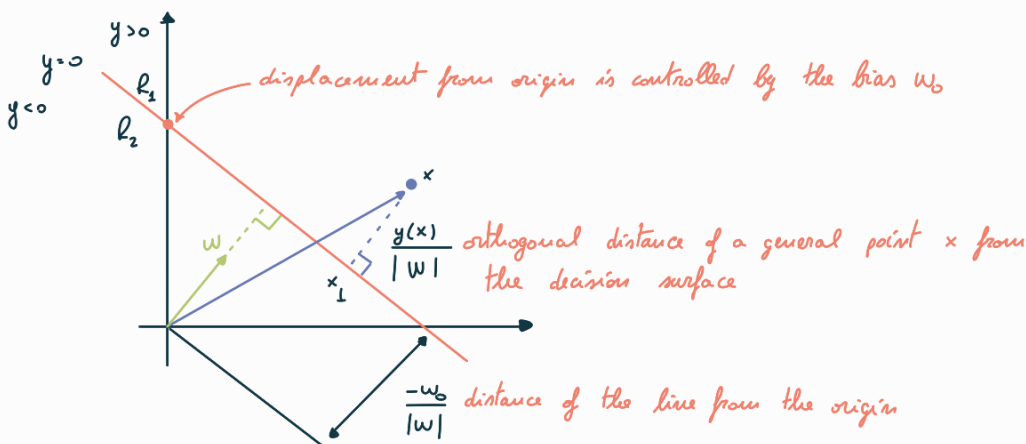
e.g.   $y(\tilde{x}, \tilde{W}) = w_0 + w_1 x + w_2 x^2 + w_3 x_3^3$
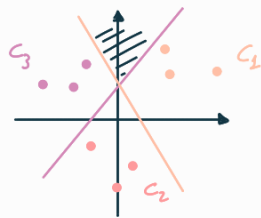
non linear $x$    linear $\tilde{w}$



# Geometric interpretation



displacement from origin is controlled by the bias $w_0$

$\dfrac{y(x)}{|w|}$ orthogonal distance of a general point $x$ from the decision surface

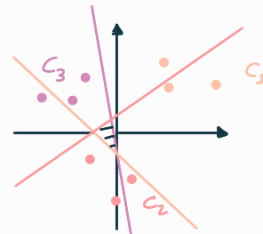$\dfrac{-w_0}{|w|}$ distance of the line from the origin

## Problem with multiple classes

we don't know how to predict the values in these regions



one vs the rest kind of classifier
$c_i$ vs not $c_i$

■ is $c_3$ or not?
■ is $c_1$ or not?

$k-1$ binary classifiers needed

one vs one kind of classifier
$c_i$ vs $c_j$

■ $c_1$ vs $c_2$ (does not care about $c_3$)
■ $c_2$ vs $c_3$ (does not care about $c_1$)
■ $c_3$ vs $c_1$ (does not care about $c_2$)

$K(k-1)/2$ binary classifiers needed

We can't reduce the problem of multi class classification to a set of binary classifiers.
We need to consider a different model suitable for K classes, such that $\longrightarrow$
Prediction is possible in every region.



## What we have

$$D = \{ (x_m, t_m)_{m=1}^{N} \}$$

one-of-k coding scheme   e.g. $t_m = (\overset{1}{0}, ..., 0, \overset{k-1}{1}, \overset{k}{0}, \overset{k+1}{...}, 0)^T \longleftrightarrow x_m \in C_k$

$t_j = 0 \; \forall j \in [1, ..., N], \; j \neq K \; ; t_k = 1$

$$y(x) = \tilde{W}^T \tilde{z} \quad \tilde{z} = \begin{pmatrix} 1 \\ x \end{pmatrix} \quad \tilde{W}^T = \begin{pmatrix} \tilde{w}_1 \\ \vdots \\ \tilde{w}_N \end{pmatrix} \quad \tilde{w}_i = \begin{pmatrix} w_{oi} \\ w_i \end{pmatrix}$$

## Approaches to learn linear k-class discriminants

Approaches:
■ Least squares
■ Perceptron
■ Linear Discriminant Analysis (LDA)
■ Support Vector Machines (SVM)

## Spoilers:

while methods such as *logistic regression* (probabilistic discriminative model) learn using the most representative samples for each class, SVMs learn using the most ambiguos and difficult to classify samples (the nearest to other classes) and use only them, ignoring the others.